

Kernel-Aware Graph Prompt Learning for Few-Shot Anomaly Detection

Fenfang Tao¹, Guo-Sen Xie^{1*}, Fang Zhao², Xiangbo Shu¹

¹School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing, China

²School of Intelligence Science and Technology, Nanjing University, Suzhou, China

{fenfangtao294, gsxiehm, zhaofang0627, shuxb104}@gmail

Abstract

Few-shot anomaly detection (FSAD) aims to detect unseen anomaly regions with the guidance of very few normal support images from the same class. Existing FSAD methods usually find anomalies by directly designing complex text prompts to align them with visual features under the prevailing large vision-language model paradigm. However, these methods, almost always, neglect intrinsic contextual information in visual features, e.g., the interaction relationships between different vision layers, which is an important clue for detecting anomalies comprehensively. To this end, we propose a kernel-aware graph prompt learning framework, termed as KAG-prompt, by reasoning the cross-layer relations among visual features for FSAD. Specifically, a kernel-aware hierarchical graph is built by taking the different layer features focusing on anomalous regions of different sizes as nodes, meanwhile, the relationships between arbitrary pairs of nodes stand for the edges of the graph. By message passing over this graph, KAG-prompt can capture cross-layer contextual information, thus leading to more accurate anomaly prediction. Moreover, to integrate the information of multiple important anomaly signals in the prediction map, we propose a novel image-level scoring method based on multi-level information fusion. Extensive experiments on MVTecAD and VisA datasets show that KAG-prompt achieves state-of-the-art FSAD results for image-level/pixel-level anomaly detection.

Code — <https://github.com/CVL-hub/KAG-prompt.git>

Introduction

Industrial anomaly detection (AD) (Roth et al. 2022), as an important task in the computer vision field, plays a crucial role in modern manufacturing and production. The AD task aims to detect and localize anomalies in industrial product images. Since real-world anomalies from different scenarios (Bergmann et al. 2019; Zou et al. 2022) usually differ in texture, color, shape, size, etc., acquiring and labeling these defects is costly and labor-intensive. Typically, AD models are directly trained on normal samples to identify anomalous samples (Zavrtanik, Kristan, and Skočaj 2021). However, this ideal setting seldom exists in real-world situations,

*Corresponding Author

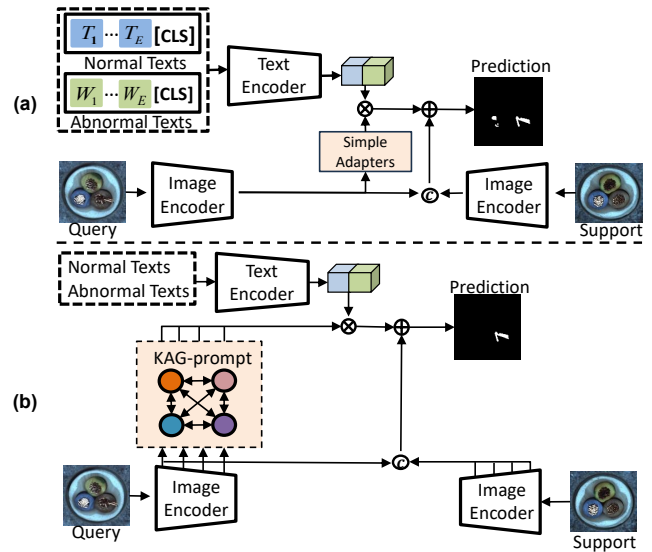


Figure 1: **Comparisons of KAG-prompt and existing FSAD models.** (a) Existing FSAD methods usually design complex text prompts, i.e., T_i, W_i are manually designed and/or learnable text prompts. For query image branches, they only learn simple adapters to extract visual features for downstream tasks. However, this paradigm segments normal backgrounds into anomaly ones. (b) Our KAG-prompt can well predict the anomalies in the query image by constructing a kernel-aware hierarchical graph to capture cross-layer multi-level relationships.

e.g., when there is rare or even no relevant training data in a new industry chain, inspecting product defects becomes infeasible by the aforementioned routine. As such, FSAD is proposed to alleviate the annotation costs and novel anomaly detection issues.

FSAD aims to detect the anomaly region in a query image by additionally utilizing a few normal support images, where the object class of the support/query images is the same. Recently, large vision-language models (LVLMs) (Zhu et al. 2023; Su et al. 2023), e.g., CLIP (Radford et al. 2021), have demonstrated robust capabilities in recognizing unknown objects (Gu et al. 2021; Kuo et al. 2022) and detecting out-

of-distribution data (Cohen, Abutbul, and Hoshen 2022), which paves the way for advancing AD and FSAD tasks (Xie et al. 2021a,b) and leads to seminal prompt learning based methods (Jeong et al. 2023; Zhou and Wang 2024; Wang et al. 2022). Most of these methods utilize a two-branch network for taking text and query image as input respectively, and the anomaly map is obtained by performing matrix operation among text and image features. For the text branch, a large number of artificial text prompts are usually designed and/or learned, and then aggregated to achieve the normal/abnormal text features (Fig. 1(a)). Since the above process is heuristic and/or partially learned, it is difficult to encompass all types of anomalies, e.g., some background region is mis-segmented as anomaly ones in Fig. 1(a). In addition, these text-prompt engineering methods involve human intervention and elaborate design, and thus usually cannot meet the automation requirements in real-world industrial scenarios. For the query image branch, existing methods merely rely on simple adapters (Chen et al. 2023; Gu et al. 2024b) to extract image features for afterward matrix operation, e.g., by using sliding windows (Jeong et al. 2023) or kernels of different sizes (Gu et al. 2024a) for detecting anomaly regions. Although, multi-layer features (Chen et al. 2023; Li et al. 2024b), e.g., PromptAD (Li et al. 2024b), have been adopted for accurate anomaly estimates. However, the high-order contextual relationships among different vision layers are still not fully leveraged for better inferring query anomaly regions.

To address the above challenges, we propose a kernel-aware graph prompt learning framework (KAG-prompt) to comprehensively reason the high-order relationships among cross-layer features for FSAD. KAG-prompt consists of a kernel-aware hierarchical graph (KAHG) module and a multi-information fusion (MIF) module, respectively. As in Fig. 1(b), solely relying on simple manually designed text prompts, KAG-prompt implicitly learns visual prompts by constructing a prompt graph among different layer features for capturing multi-level contextual information. In this manner, KAG-prompt can transfer more relation context from the seen to unseen domains. Specifically, KAHG first uses convolutional kernels of different sizes for extracting visual features of different layers. These kernel-aware features from each layer represent a node in the graph accordingly, and each node w.r.t. visual features of different layers indicate anomalous regions of different sizes. Meanwhile, the relationships between arbitrary pairs of nodes stand for the edges of the graph. Through message passing over this graph to encourage desirable information interactions among different visual layers, KAG-prompt mines richer relationships between cross-layer features in a structured way, thus leading to a more complete understanding of the anomalous region for accurate localization (in Fig.1(b), the anomaly region is well estimated in the query image). In addition, traditional methods often rely on the maximum value to assess the anomalousness of the prediction map, however, they may overlook other potentially important information. As such, we also propose to utilize the average of the top-k maxima in the prediction map as a new metric, which effectively integrates the information of multiple

important anomaly signals in the prediction map, thus improving the accuracy and robustness of anomaly detection. This top-k-induced strategy together with the global class score leads to our multi-information fusion module.

In summary, our contributions are as follows:

- We propose kernel-aware graph prompt learning (KAG-prompt), which constructs a kernel-aware hierarchical graph to capture cross-layer multi-level visual relations, thereby enhancing anomaly detection and recognition capabilities.
- We propose a novel image-level scoring method based on multi-level information fusion. By averaging the top-k maximum values from prediction maps and weighting them with scores from global feature-text alignment calculations, our approach effectively integrates critical anomaly signals from prediction maps, distinguishing it from traditional methods that solely rely on the maximum value.
- We reveal the significant potential of leveraging contextual information from different layer features for industrial anomaly detection. Diverging from previous methods that rely on complex text prompts, this inspires us to propose a novel FSAD method based on graph prompt learning.
- Extensive experiments on MVTecAD (Bergmann et al. 2019) and VisA (Zou et al. 2022) datasets demonstrate that our method achieves state-of-the-art FSAD results. We further validate the effectiveness of KAG-prompt through comprehensive ablation studies.

Related Work

Anomaly Detection. AD methods can be classified into four categories: synthesis-based methods, embedding-based methods, reconstruction-based methods, and knowledge distillation-based methods. Synthesis-based methods create anomalous images by introducing noise into the images to train the model. DRAEM (Zavrtanik, Kristan, and Skočaj 2021) generates anomalies by sampling from an external texture dataset (Cimpoi et al. 2014) and producing an anomaly mask using Perlin (Perlin 1985) noise. Cut-paste (Li et al. 2021) involves cropping parts of an image and pasting them onto different regions of the same image. NSA (Schlüter et al. 2022) employs Poisson image editing to integrate patches from various images seamlessly. Additionally, some methods (Zavrtanik, Kristan, and Skočaj 2022; Liu et al. 2023) introduce noise at the feature level to synthesize anomalies.

Embedding-based methods typically utilize a pre-trained model on ImageNet (Deng et al. 2009) to encode normal samples into a high-dimensional feature space and then compute distances in this space to detect anomalies. Padim (Defard et al. 2021) employs a multivariate Gaussian distribution to estimate the probability distribution of normal samples. Patchcore (Roth et al. 2022) stores representative features in the memory bank. CS-Flow (Rudolph et al. 2022) processes multiple features at different scales simultaneously while using normalized flow as the latent space. CFlow (Gudovskiy, Ishizaka, and Kozuka 2022), also based

on normalized flow, includes a discriminative encoder and a multi-scale generative decoder. InReaCh (McIntosh and Albu 2023) associates patches to channels, considering only channels with high span and low propagation as normal.

Reconstruction-based methods primarily focus on reconstructing normal inputs, with significant reconstruction errors observed in anomalies. Earlier methods (Dehaene and Eline 2020; Wang et al. 2020; Schlegl et al. 2017; Liang et al. 2023; Liu et al. 2022) usually employ AE (Rudolph, Wandt, and Rosenhahn 2019), VAE (Kingma and Welling 2013) and GAN (Creswell et al. 2018). Inspired by the success of diffusion models (Rombach et al. 2022; Ho, Jain, and Abbeel 2020) in generating high-quality and diverse images, more and more methods utilize diffusion models for anomaly modeling. AnoDDPM (Wyatt et al. 2022) is the first to apply diffusion models to the AD task. DiffusionAD (Zhang et al. 2023) utilizes generated anomaly samples and labels to achieve denoising and segmentation through two sub-networks. DiAD (He et al. 2024) uses SG networks to reconstruct anomalous regions while preserving the semantic information of the original image.

Knowledge distillation-based methods enable the student network to learn exclusively from the normal samples provided by the teacher network, with anomaly detection achieved by examining the discrepancies between the teacher and student models. RD (Deng and Li 2022) employs a reverse distillation paradigm where the input of the student network becomes the embedding of the teacher model, to recover the multi-scale representation of the teacher. RD++ (Tien et al. 2023) introduces a multi-scale projection layer based on RD and incorporates several loss constraints. MemKD (Gu et al. 2023) reinforces the normalcy of students’ extracted features by recalling stored normal information.

These aforementioned methods require a substantial quantity of normal samples to accurately model their distribution, rendering them unsuitable for dynamic production environments. In contrast, our FSAD approach uses only a small number of normal samples for inference.

Few-Shot Anomaly Detection. FSAD is first investigated by RegAD (Huang et al. 2022). Patchcore (Roth et al. 2022) also demonstrates FSAD performance but with poor results. FastRecon (Fang et al. 2023) employs distribution regularization to derive the optimal transformation from support image features to query image features. MuSc (Li et al. 2024a) utilizes patches of test images to evaluate each other, thereby constructing a normal distribution. WinCLIP (Jeong et al. 2023) develops the potential for language-driven anomaly detection by manually designing text prompts for both normal and anomalous cases and utilizing sliding windows to extract and aggregate multi-scale image features. APRIL-GAN (Chen et al. 2023) utilizes learnable linear layers to align patch-level image features with textual features, addressing the inefficiencies associated with WinCLIP’s multiple windows and further enhancing performance. AnomalyGPT (Gu et al. 2024b) proposes a decoder based on visual and textual feature matching to generate pixel-level anomaly localization results. The original image and decoder output serve as inputs to LVLM

for anomaly detection, eliminating the need for manual threshold setting. PromptAD (Li et al. 2024b) constructs a large number of negative samples by concatenating normal prompts with abnormal suffixes, thereby guiding text prompt learning. Additionally, it introduces the concept of explicit abnormal edges.

These FSAD methods focus on how to adjust the text prompts to align with the image, ignoring the fact that the images in the downstream task are very different from natural images, and the relations among visual features from different layers are not well leveraged. In contrast, our approach utilizes graph prompt learning to endow desirable informative interactions between different layers, capturing cross-layer multi-level visual relations so that the updated features can be better aligned with texts.

Method

Overview

An overview of our proposed KAG-prompt is shown in Fig. 2. Given a query image $x \in \mathbb{R}^{C \times H \times W}$, it is fed into the image encoder to obtain the patch features $P_i, i \in \{1, 2, 3, 4\}$ in each layer. P_i is first passed through the linear layer and then fed to the multi-kernel convolution to obtain V_i which focuses on anomalous regions of different sizes. V_i is used as an initial node to construct a kernel-aware hierarchical graph, which enables informative interactions among cross-layer visual features through a message passing mechanism to achieve updated visual feature N_i . Finally, N_i is aligned with textual features to obtain pixel-level anomaly localization result M_p . Additionally, storing patch features of each layer of normal samples in a memory bank, the anomaly localization result M_v is obtained by measuring the similarity between query image patches and the most similar patches in the memory bank. The final anomaly localization result M is the fusion of M_p and M_v . For image-level scoring, the cls token is aligned with text features after adapter to get image-level score s_1 , meanwhile, the average of the top-k maximum values of M is taken as s_2 ; finally, s_1 and s_2 are fused to get the final image-level score.

Kernel-Aware Hierarchical Graph

Multi-Kernel Convolution. Inspired by FiLo (Gu et al. 2024a), we also use multi-shape and multi-scale convolutional kernels to focus on anomalous regions of different sizes. Specifically, for a query image $x \in \mathbb{R}^{C \times H \times W}$, which is input to an image encoder, the intermediate patch feature P_i is obtained at each stage $i, i \in \{1, 2, 3, 4\}$. Since the intermediate features are not subjected to final image-text alignment, we first pass P_i through the linear layer so that it is in the same dimension as the text, followed by convolutional kernels of different shapes and sizes to obtain $V_{i,k}$. Subsequently, these features are aggregated and normalized to obtain V_i . At this point, V_i not only has the semantic information of the level but also pays attention to the anomalous regions of different shapes and sizes. V_i is defined as follows:

$$V_i = \text{Norm} \left(\sum_{k=1}^n C_k (\text{Linear}(P_i)) \right), \quad (1)$$

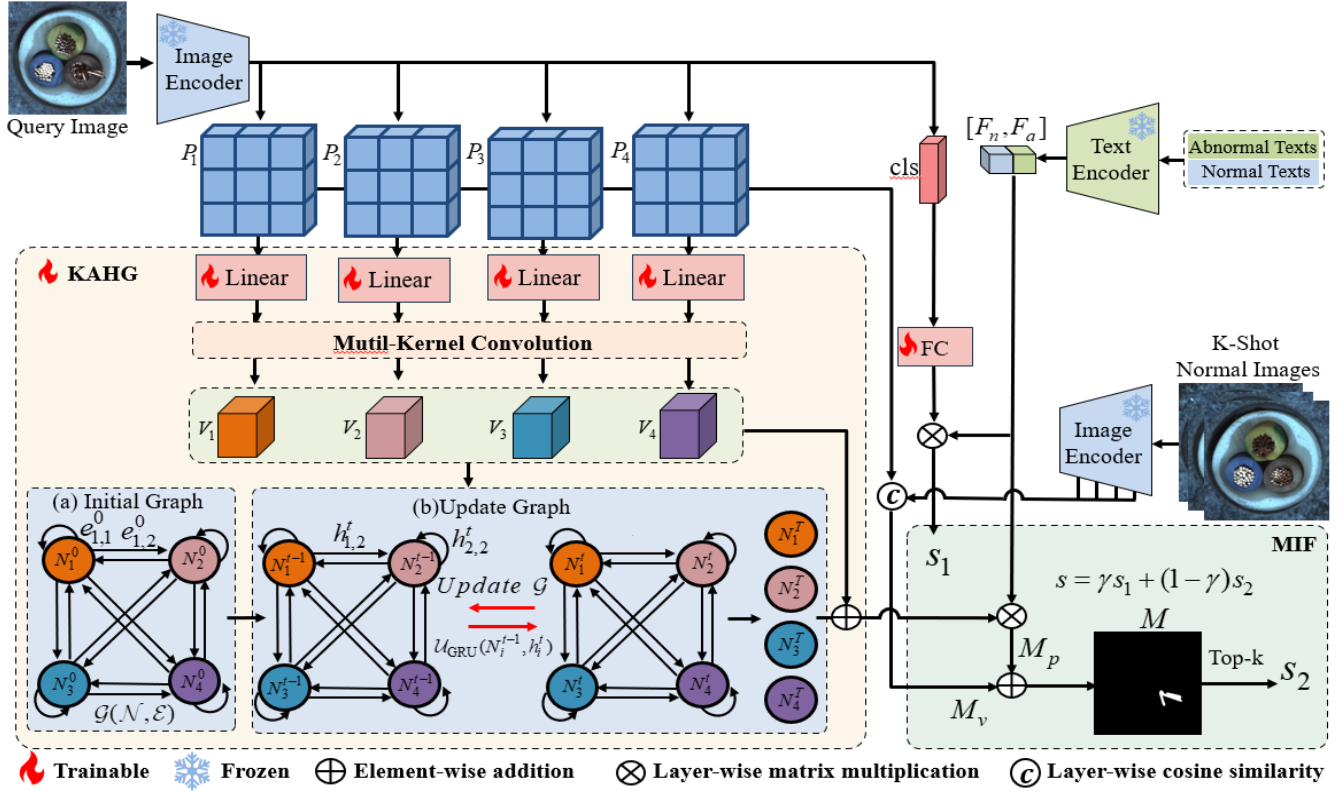


Figure 2: **The architecture of KAG-prompt.** KAG-prompt contains two modules, i.e., KAHG and MIF. The KAHG module takes visual features from different layers as input and these features undergo information interaction within the kernel-aware hierarchical graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ before aligning with texts to obtain an anomaly localization map M_p . Next, the distance between the query image and the most similar patch feature in the memory bank is calculated to get the localization map M_v . In the MIF module, for image-level score calculation, the cls token is first adapted and aligned with the texts to get s_1 ; then, M_p and M_v are fused to get s_2 by a top-k fusion mechanism; finally, s_1 and s_2 are fused to achieve the image-level score s .

where the n different convolutional kernels are denoted as $C_k(\cdot)$ and k takes values from 1 to n , and $\text{Norm}(\cdot)$ represents the normalization operation.

Kernel-Aware Hierarchical Node Embedding. We use the patch features of each layer that have gone through the multi-kernel convolution as nodes. For node N_i , its initialization N_i^0 is denoted as:

$$N_i^0 = V_i \in \mathbb{R}^{C' \times H' \times W'}, \quad (2)$$

where C' is the channel, and $H' \times W'$ is the spatial resolution.

Edge Embedding. The kernel-aware hierarchical graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ is a fully connected graph with self-connections. A loop-edge $e_{i,i}$ is an edge that connects itself. Specifically, we utilize the intra-attention mechanism to compute the response at a given position by focusing on all positions of that node to capture the internal structural relationships of the layer's patch features.

$$\begin{aligned} e_{i,i}^t &= \mathcal{F}_{\text{intra-att}}(N_i^t) \in \mathbb{R}^{C' \times H' \times W'} \\ &= \alpha \text{softmax}((W_{c_1} * N_i^t)(W_{c_2} * N_i^t)^T)(W_{c_3} * N_i^t) \\ &\quad + N_i^t, \end{aligned} \quad (3)$$

where α is a learnable scale parameter, W_{c_j} denotes the learnable convolutional kernel and $*$ denotes the convolution operation.

The line-edge $e_{i,j}$ represents the directed connection from node N_i to N_j and is used to capture the relationship between them. Specifically, we utilize the inter-attention mechanism to construct the line-edge:

$$e_{j,i}^t = \mathcal{F}_{\text{inter-att}}(N_i^t, N_j^t) = N_i^t W_c N_j^{tT} \in \mathbb{R}^{W' H' \times W' H'}, \quad (4)$$

$$e_{j,i}^t = \mathcal{F}_{\text{inter-att}}(N_j^t, N_i^t) = N_j^t W_c^T N_i^{tT} \in \mathbb{R}^{W' H' \times W' H'}, \quad (5)$$

where $W_c \in \mathbb{R}^{C' \times C'}$ is the learnable weight matrix. N_i^t and N_j^t are flattened into matrices of shape $W' H' \times C'$. By focusing on each pair of nodes, $e_{i,j}^t$ reacts to the remote relationships between nodes in different layers.

Hierarchical Node Message Passing. Since the loop-edge $e_{i,i}$ itself contains the raw and contextual information of the node in that layer, we consider itself as the message $h_{i,i}^t$ delivered by the loop-edge:

$$h_{i,i}^t = e_{i,i}^{t-1} \in \mathbb{R}^{C' \times H' \times W'}. \quad (6)$$

For the message passed from node N_j to N_i , we instead assign its edges as weighted weights to neighboring nodes:

$$\begin{aligned} h_{j,i}^t &= \mathcal{F}_{\text{mes}}(N_j^{t-1}, e_{i,j}^{t-1}) \\ &= \mathcal{F}_{\text{reshape}}(\text{softmax}(e_{i,j}^{t-1})N_j^{t-1}) \in \mathbb{R}^{C' \times H' \times W'}, \end{aligned} \quad (7)$$

where $\text{softmax}(\cdot)$ is normalized to each row of the input.

Considering that the noise present in layer-based nodes can adversely affect the information in the message passing, we use learnable gating to measure the confidence of the message:

$$a_{j,i}^t = \mathcal{F}_{\text{gate}}(h_{j,i}^t) = \sigma(\mathcal{F}_{\text{GAP}}(W_{\text{GAP}} * h_{j,i}^t + b)) \in [0, 1], \quad (8)$$

where $\mathcal{F}_{\text{GAP}}(\cdot)$ is the global average pooling operation, W_{GAP} and b are its convolutional kernel and bias, and σ is the sigmoid activation function.

Node N_i^t receives the total messages from neighboring nodes and itself through gating as follows:

$$h_i^t = \sum_{j \in \mathcal{N}} a_{j,i}^t \odot h_{j,i}^t \in \mathbb{R}^{C' \times H' \times W'}, \quad (9)$$

where \odot is the channel-by-element product.

Hierarchical Nodes Updates. In iteration t , we utilize ConvGRU (Ballas et al. 2015), which updates the node state by aggregating the obtained total messages (Eq. (9)) as well as the node state at step $t - 1$:

$$N_i^t = \mathcal{U}_{\text{GRU}}(N_i^{t-1}, h_i^t) \in \mathbb{R}^{C' \times H' \times W'}. \quad (10)$$

After the message has been iterated T times, we output the node N_i^T . To preserve the hierarchical and contextual information of the original input, we residually concatenate N_i^T with the original input V_i (Eq. (1)) to obtain O_i :

$$O_i = \mathcal{F}_{\text{reshape}}(N_i^T + V_i) \in \mathbb{R}^{W' \times H' \times C'}. \quad (11)$$

Anomaly Detection By Multi-Information Fusion

We use manually designed textual prompts (Jeong et al. 2023), processed through a text encoder to obtain normal and abnormal text features, denoted as $F_{\text{text}} = [F_n, F_a] \in \mathbb{R}^{2 \times C'}$. We align the layers of visual features with rich contextual information to the text prompts, then aggregate and normalize the predictions M_{p_i} from each layer to generate anomaly map M_p .

$$M_p = \text{Up}(\text{Norm} \sum_{i=1}^4 \text{softmax}(O_i F_a^T)), \quad (12)$$

where $\text{Up}(\cdot)$ denotes Bilinear interpolation.

Additionally, we store the visual features of the support set across layers in a memory bank, denoted as R . The localization result M_v is obtained by calculating the distance between the query patch and the most similar corresponding patch in R .

$$M_v = \text{Up}(\sum_{i=1}^4 (1 - \max_{r \in R} (O_i r_i))). \quad (13)$$

The final result of the anomaly localization is:

$$M = \gamma M_p + (1 - \gamma) M_v. \quad (14)$$

For image-level scoring, unlike the previous reliance on maxima to assess the anomalousness of the prediction map, this may overlook other potentially important information. For this reason, we propose utilizing the average of the top-k maximum values in the prediction map as a new metric, the approach that effectively synthesizes information from multiple important anomaly signals in the prediction map. Specifically, the cls feature is aligned to the text features after passing through an adapter (fully connected layer, FC) to obtain a global score s_1 . Next, we average the top-k maxima in M (Eq. (14)) to obtain s_2 . The final image-level score is the fusion of s_1 and s_2 :

$$s_1 = \text{FC}(F_{\text{cls}}) F_a^T, \quad (15)$$

$$s = \gamma s_1 + (1 - \gamma) \text{Mean}(\text{Top-k}(M)). \quad (16)$$

Loss Function

We mainly use cross-entropy loss, focal loss (Lin et al. 2017), and dice loss (Milletari, Navab, and Ahmadi 2016). The overall losses are as follows:

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{\text{CE}}(s_1, c) + \lambda_1 \sum_{i=1}^4 (\mathcal{L}_{\text{Focal}}([I - M_{p_i}, M_{p_i}], G)) \\ &\quad + \lambda_2 \sum_{i=1}^4 (\mathcal{L}_{\text{Dice}}(M_{p_i}, G) + \mathcal{L}_{\text{Dice}}(I - M_{p_i}, I - G)), \end{aligned} \quad (17)$$

where λ_1 and λ_2 are set to 1.0 in all the experiments, $[\cdot, \cdot]$ denotes concatenation along the channels, G represents ground truth, c denotes true label, and I represents an all-ones matrix.

Experiments

Settings

Datasets. We mainly conduct experiments on the MVTEC-AD (Bergmann et al. 2019) and VisA (Zou et al. 2022) datasets. The MVTEC-AD dataset contains 5,354 high-resolution images of 5 textures and 10 objects. The training set contains 3,629 sample images without anomalies. The test set contains 1,725 images including normal and anomalous samples. The VisA dataset has 12 subsets containing 10,821 high-resolution images, of which 9,621 are normal images and 1,200 are anomalous images. As with AnomalyGPT (Gu et al. 2024b), we use the training set of one dataset as well as the synthesized anomalous images for training and perform few-shot testing on the other dataset.

Evaluation Metrics. We use area under the receiver operating characteristic (AUROC) as an image-level anomaly detection metric. In addition, we use pixel-wise AUROC (pAUROC) to evaluate anomaly localization.

Implementation Details. Our baseline model is AnomalyGPT (Gu et al. 2024b). We synthesize anomaly data for

Setup	Method	Public	MVTecAD		VisA		avg
			AUROC	pAUROC	AUROC	pAUROC	
1-shot	SPADE (Cohen et al. 2020)	arXiv2020	81.0	91.2	79.5	95.6	86.8
	PatchCore (Roth et al. 2022)	CVPR2022	83.4	92.0	79.9	95.4	87.7
	WinCLIP (Jeong et al. 2023)	CVPR2023	93.1	95.2	83.8	96.4	92.1
	APRIL-GAN (Chen et al. 2023)	arXiv2023	92.0	95.1	<u>91.2</u>	96.0	<u>93.6</u>
	AnomalyGPT† (Gu et al. 2024b)	AAAI2024	94.1	95.3	87.4	96.2	93.3
	PromptAD (Li et al. 2024b)	CVPR2024	<u>94.6</u>	<u>95.9</u>	86.9	<u>96.7</u>	93.5
	KAG-prompt (ours)	-	95.8	96.2	91.6	97.0	95.2
2-shot	SPADE (Cohen et al. 2020)	arXiv2020	82.9	92.0	80.7	96.2	88.0
	PatchCore (Roth et al. 2022)	CVPR2022	86.3	93.3	81.6	96.1	89.3
	WinCLIP (Jeong et al. 2023)	CVPR2023	94.4	96.0	84.6	96.8	93.0
	APRIL-GAN (Chen et al. 2023)	arXiv2023	92.4	95.0	<u>92.2</u>	96.2	94.0
	AnomalyGPT† (Gu et al. 2024b)	AAAI2024	95.5	95.6	88.6	96.4	94.0
	PromptAD (Li et al. 2024b)	CVPR2024	<u>95.7</u>	<u>96.2</u>	88.3	<u>97.1</u>	<u>94.3</u>
	KAG-prompt (ours)	-	96.6	96.5	92.7	97.4	95.8
4-shot	SPADE (Cohen et al. 2020)	arXiv2020	84.8	92.7	81.7	96.6	89.0
	PatchCore (Roth et al. 2022)	CVPR2022	88.8	94.3	85.3	96.8	91.3
	WinCLIP (Jeong et al. 2023)	CVPR2023	95.2	96.2	87.3	97.2	94.0
	APRIL-GAN (Chen et al. 2023)	arXiv2023	92.8	95.9	<u>92.2</u>	96.2	94.3
	AnomalyGPT† (Gu et al. 2024b)	AAAI2024	96.3	96.2	90.6	96.7	<u>95.0</u>
	PromptAD (Li et al. 2024b)	CVPR2024	<u>96.6</u>	<u>96.5</u>	89.1	<u>97.4</u>	94.9
	KAG-prompt (ours)	-	97.1	96.7	93.3	97.7	96.2

Table 1: Performance comparisons of the FSAD methods on the MVTecAD and VisA datasets. Bold indicates the best performance and underlining indicates sub-optimal results. † indicates our baseline.

KAHG	Global s_1	Max	Top-k Fusion	AUROC	pAUROC
		✓		89.7	96.5
✓		✓		91.3	97.0
✓	✓	✓		91.4	97.0
✓	✓		✓	91.6	97.0

Table 2: Module ablation at the 1-shot setting of the VisA dataset.

each normal image by NSA (Schlüter et al. 2022) technique and use it to train the model. The image resolution is 224×224. We extract visual features from layers 8, 16, 24, and 36 of the image encoder ImageBind-Huge (Girdhar et al. 2023). During training, the learning rate is set to 1e-3, batch size to 16, and the number of iterations T for graph prompts is 5. Two RTX-3090 GPUs are used for acceleration during training. The model is trained for 50 epochs on the MVTecAD dataset and 80 epochs on the VisA dataset. We set the fusion coefficient γ to 0.1 and select the top-30 scores using a top-k strategy.

Comparisons with State-of-the-Arts

Tab. ?? demonstrates the comparison results of KAG-prompt with existing few-shot anomaly detection methods SPADE, PatchCore, Winclip, AnomalyGPT, APRIL-GAN

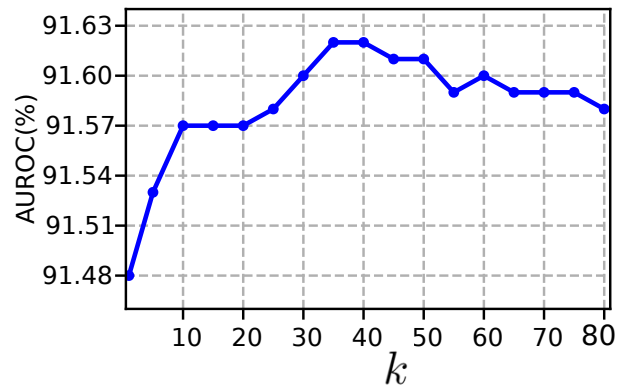


Figure 3: Ablation on top-k strategy k on the 1-shot setting of the VisA dataset.

and PromptAD. KAG-prompt shows significant improvements over all methods across both datasets in all metrics. Notably, we achieve a 1.2% improvement in AUROC and a 0.3% improvement in pAUROC on MVTecAD with the 1-shot setting compared to the suboptimal method PromptAD. Similarly, KAG-prompt achieves a 4.7% improvement in AUROC and a 0.3% improvement in pAUROC compared to PromptAD in the 1-shot setting of VisA. Considering that PromptAD has a higher pixel-level metric but a lower image-level metric on VisA, APRIL-GAN shows better per-

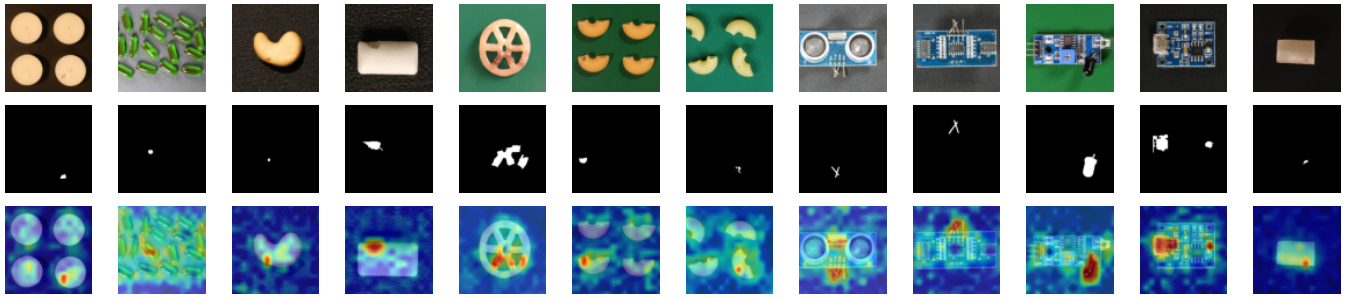


Figure 4: Visualization of KAG-prompt on VisA under 1-shot setting. The first row shows the query image, the second row depicts the corresponding ground truth, and the third row displays the heatmap of abnormal localization by KAG-prompt.

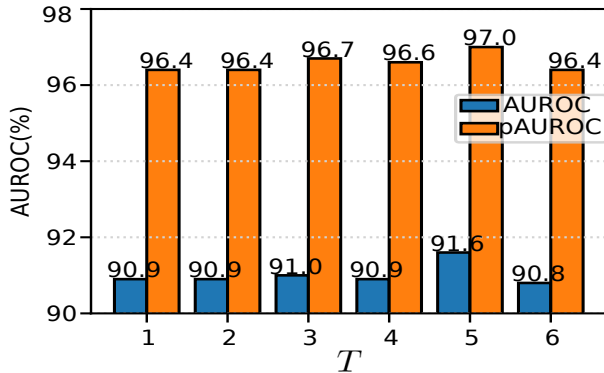


Figure 5: Ablation on graph prompt iterations T at the 1-shot setting on the VisA dataset.

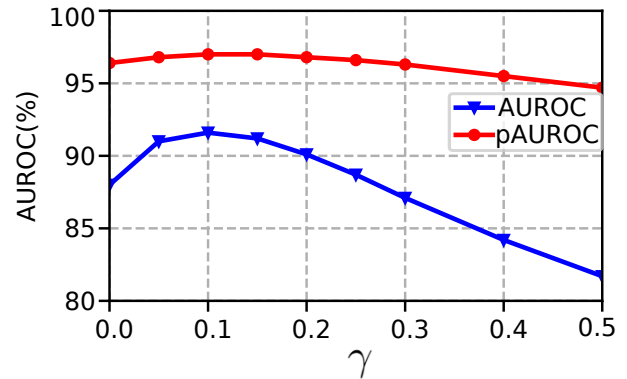


Figure 6: Ablation on the fusion coefficient γ at the 1-shot setting on the VisA dataset.

formance in both metrics. Compared to APRIL-GAN, our method improves AUROC by 0.4% and pAUROC by 1% on VisA. After averaging the results of all the metrics, our method shows the best results. Compared to the sub-optimal results, KAG-prompt is better by 1.6%. KAG-prompt similarly reaches the state-of-the-art in 2-/4-shot settings. The quantitative results of KAG-prompt in anomaly localization demonstrate its significant advantages in terms of performance under different numbers of support images.

Ablation Study

Module Ablation. We first verify the effectiveness of different modules of KAG-prompt, including the baseline maximum value calculation (Max), kernel-aware hierarchical graph (KAHG), the cls-guided global scoring s_1 (Global s_1), and the top-k strategy for multi-information fusion (Top-k Fusion). The results are shown in Tab. ??, where each module contributes to the superior performance of KAG-prompt, with kernel-aware hierarchical graph being the most important one. Compared to the baseline, it improves AUROC by 1.6% and pAUROC by 0.5%.

Graph Prompt Iterations T . We vary T from 1 to 6 to observe the performance of KAG-prompt. As in Fig. 5, $T = 5$ performs best, thus we set $T = 5$ in all our experiments.

Top-k Strategy k . Fig. 3 illustrates the impact of varying values of k on inference performance. We vary k from 1 to 80 and obtain the best result 91.62% when $k = 40$, after

which the performance decreases as k increases. As such, we set $k = 30$ in all experiments.

The Fusion Coefficient γ . We vary γ from 0 to 0.5, as shown in Fig. 6. The results are best at $\gamma = 0.1$. Therefore, we set $\gamma = 0.1$ for all experiments.

Visualization Results

Fig. 4 shows the visualization results of KAG-prompt on the VisA dataset. KAG-prompt has achieved anomaly localization results that closely match their ground truth. This demonstrates the exceptional anomaly localization capability of KAG-prompt.

Conclusion

In this paper, we propose a novel anomaly detection method, KAG-prompt. KAG-prompt constructs a kernel-aware hierarchical graph, which learns contextual relationships between cross-layer hierarchical visual features while focusing on anomalous regions of different sizes, to extract updated visual features for aligning them with texts. In addition, vision-guided anomaly detection is introduced to improve the accuracy and robustness of anomaly detection by integrating the information of multiple important anomaly signals through multi-information fusion. Experiments on two commonly used datasets demonstrate the effectiveness of KAG-prompt under FSAD setting.

Acknowledgments

This work is supported by the National Natural Science Foundation of China (Nos. 62276134, 62476124, and 62072245).

References

- Ballas, N.; Yao, L.; Pal, C.; and Courville, A. 2015. Delving deeper into convolutional networks for learning video representations. *arXiv preprint arXiv:1511.06432*.
- Bergmann, P.; Fauser, M.; Sattlegger, D.; and Steger, C. 2019. MVTec AD—A comprehensive real-world dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9592–9600.
- Chen; et al. 2023. A zero-/fewshot anomaly classification and segmentation method for cvpr 2023 vand workshop challenge tracks 1&2: 1st place on zero-shot ad and 4th place on few-shot ad. *arXiv preprint arXiv:2305.17382*, 2(4).
- Cimpoi, M.; Maji, S.; Kokkinos, I.; Mohamed, S.; and Vedaldi, A. 2014. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3606–3613.
- Cohen; et al. 2020. Sub-image anomaly detection with deep pyramid correspondences. *arXiv preprint arXiv:2005.02357*.
- Cohen, N.; Abutbul, R.; and Hoshen, Y. 2022. Out-of-distribution detection without class labels. In *European Conference on Computer Vision*, 101–117. Springer.
- Creswell, A.; White, T.; Dumoulin, V.; Arulkumaran, K.; Sengupta, B.; and Bharath, A. A. 2018. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1): 53–65.
- Defard, T.; Setkov, A.; Loesch, A.; and Audigier, R. 2021. Padim: a patch distribution modeling framework for anomaly detection and localization. In *International Conference on Pattern Recognition*, 475–489. Springer.
- Dehaene, D.; and Eline, P. 2020. Anomaly localization by modeling perceptual features. *arXiv preprint arXiv:2008.05369*.
- Deng, H.; and Li, X. 2022. Anomaly detection via reverse distillation from one-class embedding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9737–9746.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Fang, Z.; Wang, X.; Li, H.; Liu, J.; Hu, Q.; and Xiao, J. 2023. Fastrecon: Few-shot industrial anomaly detection via fast feature reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 17481–17490.
- Girdhar, R.; El-Nouby, A.; Liu, Z.; Singh, M.; Alwala, K. V.; Joulin, A.; and Misra, I. 2023. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15180–15190.
- Gu, X.; Lin, T.-Y.; Kuo, W.; and Cui, Y. 2021. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*.
- Gu, Z.; Liu, L.; Chen, X.; Yi, R.; Zhang, J.; Wang, Y.; Wang, C.; Shu, A.; Jiang, G.; and Ma, L. 2023. Remembering normality: Memory-guided knowledge distillation for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 16401–16409.
- Gu, Z.; Zhu, B.; Zhu, G.; Chen, Y.; Li, H.; Tang, M.; and Wang, J. 2024a. FiLo: Zero-Shot Anomaly Detection by Fine-Grained Description and High-Quality Localization. *arXiv preprint arXiv:2404.13671*.
- Gu, Z.; Zhu, B.; Zhu, G.; Chen, Y.; Tang, M.; and Wang, J. 2024b. Anomalygpt: Detecting industrial anomalies using large vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 1932–1940.
- Gudovskiy, D.; Ishizaka, S.; and Kozuka, K. 2022. Cflowad: Real-time unsupervised anomaly detection with localization via conditional normalizing flows. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 98–107.
- He, H.; Zhang, J.; Chen, H.; Chen, X.; Li, Z.; Chen, X.; Wang, Y.; Wang, C.; and Xie, L. 2024. A diffusion-based framework for multi-class anomaly detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 8472–8480.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.
- Huang, C.; Guan, H.; Jiang, A.; Zhang, Y.; Spratling, M.; and Wang, Y.-F. 2022. Registration based few-shot anomaly detection. In *European Conference on Computer Vision*, 303–319. Springer.
- Jeong, J.; Zou, Y.; Kim, T.; Zhang, D.; Ravichandran, A.; and Dabeer, O. 2023. Winclip: Zero-/few-shot anomaly classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 19606–19616.
- Kingma, D. P.; and Welling, M. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Kuo, W.; Cui, Y.; Gu, X.; Piergiovanni, A.; and Angelova, A. 2022. F-vm: Open-vocabulary object detection upon frozen vision and language models. *arXiv preprint arXiv:2209.15639*.
- Li, C.-L.; Sohn, K.; Yoon, J.; and Pfister, T. 2021. Cutpaste: Self-supervised learning for anomaly detection and localization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9664–9674.
- Li, X.; Huang, Z.; Xue, F.; and Zhou, Y. 2024a. Musc: Zero-shot industrial anomaly classification and segmentation with mutual scoring of the unlabeled images. *arXiv preprint arXiv:2401.16753*.

- Li, X.; Zhang, Z.; Tan, X.; Chen, C.; Qu, Y.; Xie, Y.; and Ma, L. 2024b. Promptad: Learning prompts with only normal samples for few-shot anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16838–16848.
- Liang, Y.; Zhang, J.; Zhao, S.; Wu, R.; Liu, Y.; and Pan, S. 2023. Omni-frequency channel-selection representations for unsupervised anomaly detection. *IEEE Transactions on Image Processing*.
- Lin, T.-Y.; Goyal, P.; Girshick, R.; He, K.; and Dollár, P. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, 2980–2988.
- Liu, T.; Li, B.; Zhao, Z.; Du, X.; Jiang, B.; and Geng, L. 2022. Reconstruction from edge image combined with color and gradient difference for industrial surface anomaly detection. *arXiv preprint arXiv:2210.14485*.
- Liu, Z.; Zhou, Y.; Xu, Y.; and Wang, Z. 2023. Simplenet: A simple network for image anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 20402–20411.
- McIntosh, D.; and Albu, A. B. 2023. Inter-realization channels: Unsupervised anomaly detection beyond one-class classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6285–6295.
- Milletari, F.; Navab, N.; and Ahmadi, S.-A. 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, 565–571. Ieee.
- Perlin, K. 1985. An image synthesizer. *ACM Siggraph Computer Graphics*, 19(3): 287–296.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.
- Roth, K.; Pemula, L.; Zepeda, J.; Schölkopf, B.; Brox, T.; and Gehler, P. 2022. Towards total recall in industrial anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 14318–14328.
- Rudolph, M.; Wandt, B.; and Rosenhahn, B. 2019. Structuring autoencoders. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 0–0.
- Rudolph, M.; Wehrbein, T.; Rosenhahn, B.; and Wandt, B. 2022. Fully convolutional cross-scale-flows for image-based defect detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1088–1097.
- Schlegl, T.; Seeböck, P.; Waldstein, S. M.; Schmidt-Erfurth, U.; and Langs, G. 2017. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *International conference on information processing in medical imaging*, 146–157. Springer.
- Schlüter, H. M.; Tan, J.; Hou, B.; and Kainz, B. 2022. Natural synthetic anomalies for self-supervised anomaly detection and localization. In *European Conference on Computer Vision*, 474–489. Springer.
- Su, Y.; Lan, T.; Li, H.; Xu, J.; Wang, Y.; and Cai, D. 2023. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*.
- Tien, T. D.; Nguyen, A. T.; Tran, N. H.; Huy, T. D.; Duong, S.; Nguyen, C. D. T.; and Truong, S. Q. 2023. Revisiting reverse distillation for anomaly detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 24511–24520.
- Wang, L.; Zhang, D.; Guo, J.; and Han, Y. 2020. Image anomaly detection using normal data only by latent space resampling. *Applied Sciences*, 10(23): 8660.
- Wang, W.; Han, C.; Zhou, T.; and Liu, D. 2022. Visual recognition with deep nearest centroids. *arXiv preprint arXiv:2209.07383*.
- Wyatt, J.; Leach, A.; Schmon, S. M.; and Willcocks, C. G. 2022. Anoddp: Anomaly detection with denoising diffusion probabilistic models using simplex noise. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 650–656.
- Xie, G.-S.; Liu, J.; Xiong, H.; and Shao, L. 2021a. Scale-aware graph neural network for few-shot semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5475–5484.
- Xie, G.-S.; Xiong, H.; Liu, J.; Yao, Y.; and Shao, L. 2021b. Few-shot semantic segmentation with cyclic memory network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7293–7302.
- Zavrtanik, V.; Kristan, M.; and Skočaj, D. 2021. Draem-a discriminatively trained reconstruction embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, 8330–8339.
- Zavrtanik, V.; Kristan, M.; and Skočaj, D. 2022. Dsr—a dual subspace re-projection network for surface anomaly detection. In *European conference on computer vision*, 539–554. Springer.
- Zhang, H.; Wang, Z.; Wu, Z.; and Jiang, Y.-G. 2023. Diffusionad: Denoising diffusion for anomaly detection. *arXiv preprint arXiv:2303.08730*, 4.
- Zhou, T.; and Wang, W. 2024. Cross-image pixel contrasting for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Zhu, D.; Chen, J.; Shen, X.; Li, X.; and Elhoseiny, M. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.
- Zou, Y.; Jeong, J.; Pemula, L.; Zhang, D.; and Dabeer, O. 2022. Spot-the-difference self-supervised pre-training for anomaly detection and segmentation. In *European Conference on Computer Vision*, 392–408. Springer.